

Machine Learning Definition

Machine Learning is a field of AI that enables computers to learn from data **without being explicitly programmed**.

Supervised Learning Explained

Supervised learning is a type of ML where the model learns from **labeled data**.

Each training example includes an **input** and the correct **output** (label), and the model learns to map inputs to outputs.

Example:

- Input: "This is spam"
- Output (label): Spam
- Goal: Learn to classify new emails as spam or not spam.

Regression (Predicting a number)

Definition:

Regression is a type of supervised learning where the goal is to **predict a continuous value** (a number).

Example (from the image):

Problem: Predicting house price based on its size.

- **Input (X):** House size in square feet
- **Output (Y):** Price in \$1000s
- e.g., A house of size **750 ft²** → Predict price (say, **\$150K or \$200K**)

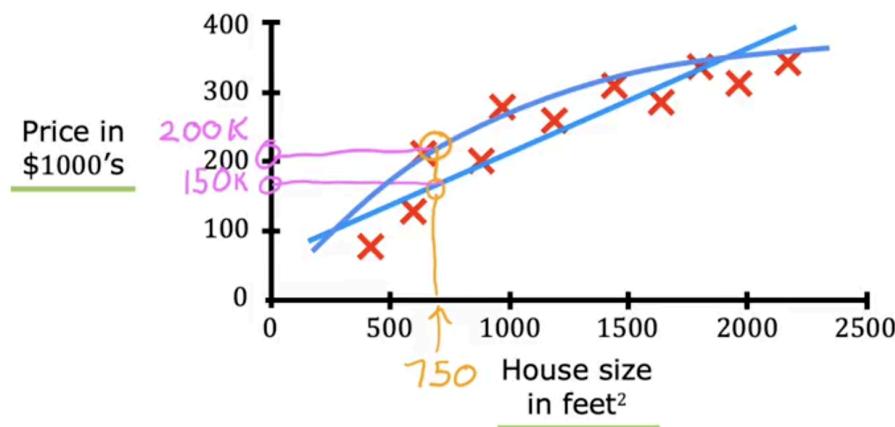
Visual Understanding:

- Red crosses = actual data points
- Blue curves = prediction models
- The model learns the pattern to estimate price for any given size.

⊕ Real-world Regression Examples:

- Predicting temperature tomorrow
- Estimating someone's weight based on height
- Predicting exam scores based on study hours

Regression: Housing price prediction



● Classification (Predicting a category)

Definition:

Classification is a type of supervised learning where the goal is to **predict a discrete label/class**, like categories (0, 1, 2...).

💡 Example: Breast Cancer Prediction

Problem: Predict whether a tumor is malignant or benign.

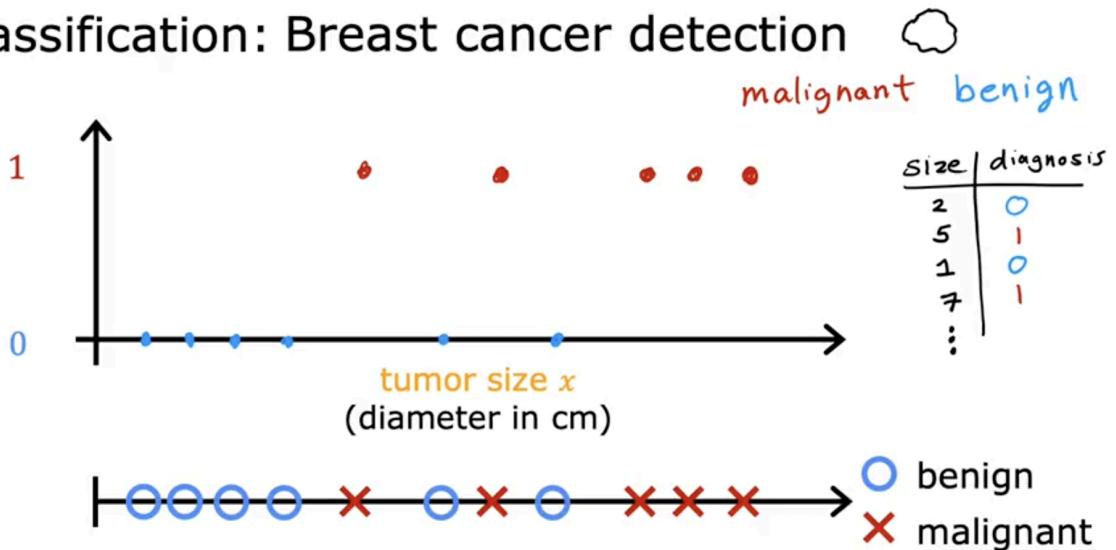
- **Input (X):** Tumor size, shape, etc. (e.g., horizontal size)
- **Output (Y):**
 - 0 → Benign (non-cancerous)
 - 1 → Malignant (cancerous)

- You never predict 0.5 or 0.75 , only **finite classes**.

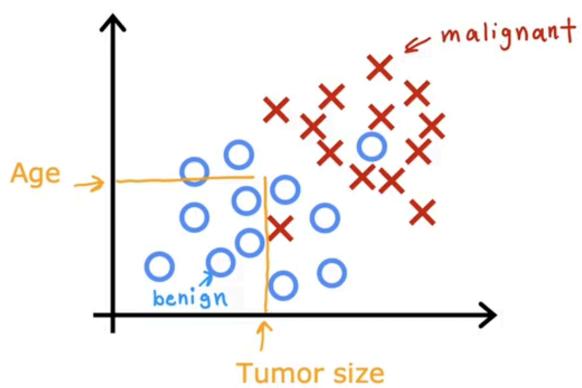
⊕ Real-world Classification Examples:

- Email: spam or not spam
- Self-driving car: identify pedestrian, car, traffic light
- Speech recognition: classify spoken word
- Visual inspection: defect or no defect in factory line

Classification: Breast cancer detection



Two or more inputs





Difference: Regression vs Classification

Feature	Regression	Classification
Output type	Continuous number (e.g., 123.45)	Discrete class (e.g., 0, 1, or 2)
Example	House price prediction	Tumor type prediction
Prediction value	Any number (e.g., 150.5K)	Limited labels (e.g., 0 or 1)
Common algorithms	Linear Regression, SVR	Logistic Regression, Decision Trees
Visualization	Line/curve fit through data points	Class boundary separation



Unsupervised Learning

Definition:

Unsupervised learning is a type of machine learning where the algorithm is given **only input data (X)** and **no labeled outputs (Y)**.

The algorithm tries to **find patterns, groupings, or structure** in the data **without supervision**.

◆ Key Concepts in Unsupervised Learning

✓ 1. Clustering → *Grouping similar data points together*

Example 1: Google News Article Clustering

- Articles about "pandas" and "zoos" appear in 5 different news pieces.
- The algorithm **detects similarities** in keywords and **groups** them into a cluster about wildlife or animals.
- No human labels are given – it **figures it out on its own**.

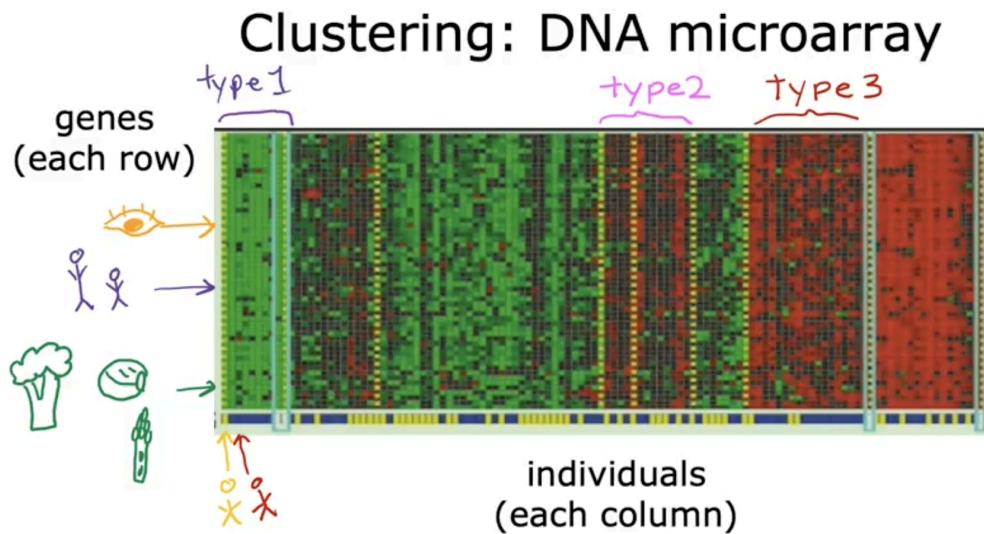


Use case: Google News automatically groups similar stories under one headline.

Example 2: DNA, Eye Color, Height

- Input data: genetic features like DNA, eye color, height, etc.
- The algorithm **groups people** based on these traits.
- It may discover "**Type 1**" humans who have similar genetic traits, **without knowing beforehand** what those types are.

📌 **Use case:** Genetic research, personalized medicine.



Example 3: Customer Segmentation in Marketing

- Input data: customer behavior like purchase history, time spent, age, income.
- The algorithm groups customers into **different market segments** like:
 - Type A: wants to grow skills
 - Type B: wants to earn money
 - Type C: stays updated
- This helps businesses **target each group separately**.

📌 **Use case:** Amazon or Netflix personalizing recommendations.

No Labels (Y):

Just raw data (X).

Model learns **structure or hidden patterns** on its own.

◆ Other Unsupervised Techniques

Anomaly Detection

- Find rare or unusual data points (outliers)
- Example: Detecting credit card fraud or server attacks

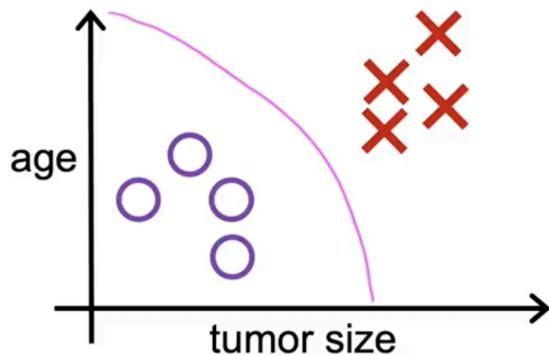
Dimensionality Reduction

- Reduce large number of features to few important ones
 - Helps visualize high-dimensional data
 - Techniques: PCA (Principal Component Analysis), t-SNE
-

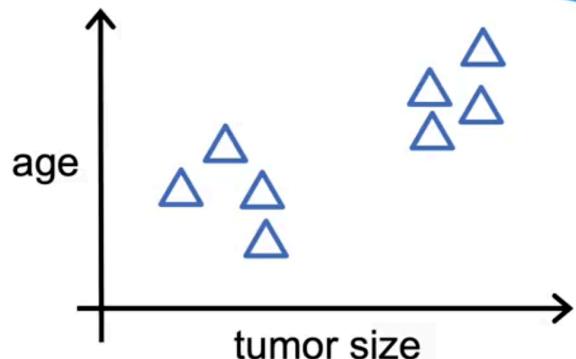
Summary Table

Feature	Unsupervised Learning
Labels provided?	 No (only input X)
Goal	Find hidden patterns/structures
Key techniques	Clustering, Anomaly Detection, PCA
Common algorithms	K-Means, DBSCAN, PCA
Real-world examples	News grouping, DNA clustering, Market segmentation

Supervised learning
Learn from data **labeled** with the “**right answers**”



Unsupervised learning
Find something interesting in **unlabeled** data.



Question

Of the following examples, which would you address using an unsupervised learning algorithm? (Check all that apply.)

- Given a set of news articles found on the web, group them into sets of articles about the same stories.

Correct

This is a type of unsupervised learning called clustering

- Given email labeled as spam/not spam, learn a spam filter.

- Given a database of customer data, automatically discover market segments and group customers into different market segments.

Correct

This is a type of unsupervised learning called clustering

- Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not

Linear Regression

- It's a **supervised learning** algorithm used for **regression tasks**, where the goal is to predict a continuous value (e.g., price).
- In this example:
Input (x) = size of the house
Output (y) = price of the house
The function $f(x) = wx + b$ ($f(x) = wx + b$) estimates price from size.

Terminology Explained

Term	Meaning
Training Set	Collection of data with inputs (features like size) and outputs (targets like price).
Learning Algorithm	The method (e.g., gradient descent) that learns the function $f(x)$ from data.
x	Input feature (e.g., house size).
$f(x)$	The function learned (also called the model or hypothesis , though "hypothesis" is crossed out in image).
\hat{y} (y-hat)	The prediction made by the model (i.e., estimated price).
y	Actual target value (i.e., real price of the house).

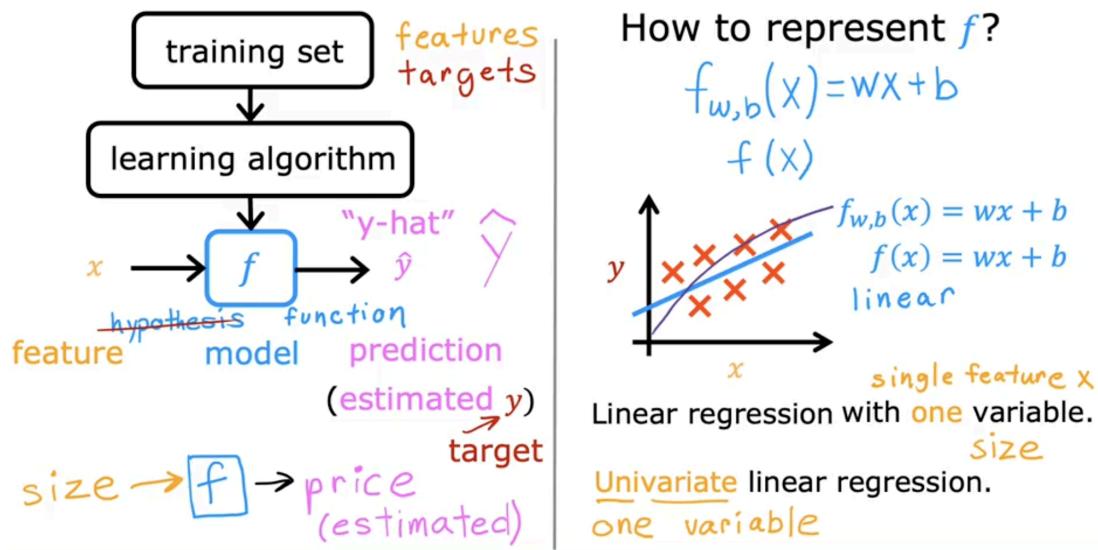
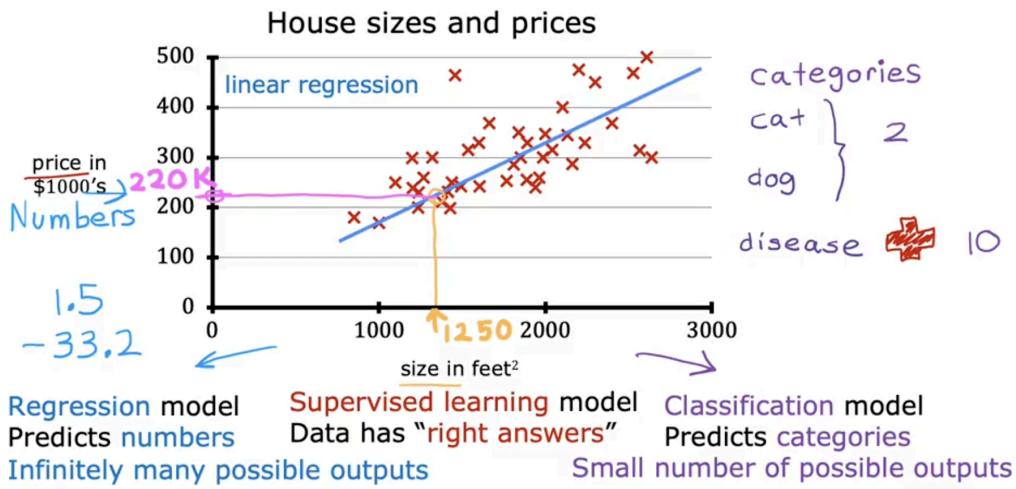
Model Representation

- **Linear function:**

$$f_{w,b}(x) = w_1x + b \quad f_{w,b}(x) = wx + b$$

where:

- w = weight (slope),
 - b = bias (intercept),
 - x = input feature (house size),
 - $f(x)$ = output prediction (house price).
-



◆ Key Points from the Image

- **Univariate Linear Regression:** Only one input variable (house size).
- **Supervised Learning:** Because we learn from labeled data (x, y).
- Graph shows red X's as real data points and a blue/purple line as the predicted function.

🎯 What is a Cost Function?

We want our model to **predict values close to real values**.

So we define a **Cost Function** (also called **Loss Function**) to measure **how wrong the predictions are**.

The goal is:

Find w and b that make our predictions close to the actual data

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

Let's break this down ↴

Symbol	Meaning
m	Number of training examples
$x^{(i)}$	Input (feature) of the i-th example
$y^{(i)}$	Actual output of the i-th example
$\hat{y}^{(i)} = f_w, b(x^{(i)}) = wx^{(i)} + b$	Predicted output
$(\hat{y}^{(i)} - y^{(i)})^2$	Squared Error for one example
$\sum_{i=1}^m$	Add all errors from all training examples
$\frac{1}{2m}$	Take the average error (divide by number of examples), and the $1/2$ is for math convenience when taking derivatives

📌 What does this function actually do?

It calculates the **average squared difference between predicted and actual outputs**.

Smaller cost ⇒ better model.

🎯 GOAL:

Find the values of w and b that make the cost function $J(w, b)$ as small as possible.

Example:

Let's say:

x	y (real)	\hat{y} (predicted)
1	2	1.5
2	4	3.5

Errors:

- For $x=1$: $(1.5 - 2)^2 = 0.25$
- For $x=2$: $(3.5 - 4)^2 = 0.25$

Cost:

$$J = \frac{1}{2 \cdot 2} (0.25 + 0.25) = \frac{0.5}{4} = 0.125$$

V1:- Linear Regression: Introduction to the Cost Function

◆ Objective of Linear Regression

- Given a **training set** with:
 - Input feature: x
 - Output target: y
- The model we use:
 $f_{w,b}(x) = w \cdot x + b$
 - w : weight (slope of the line)
 - b : bias (y-intercept)

◆ Terminology

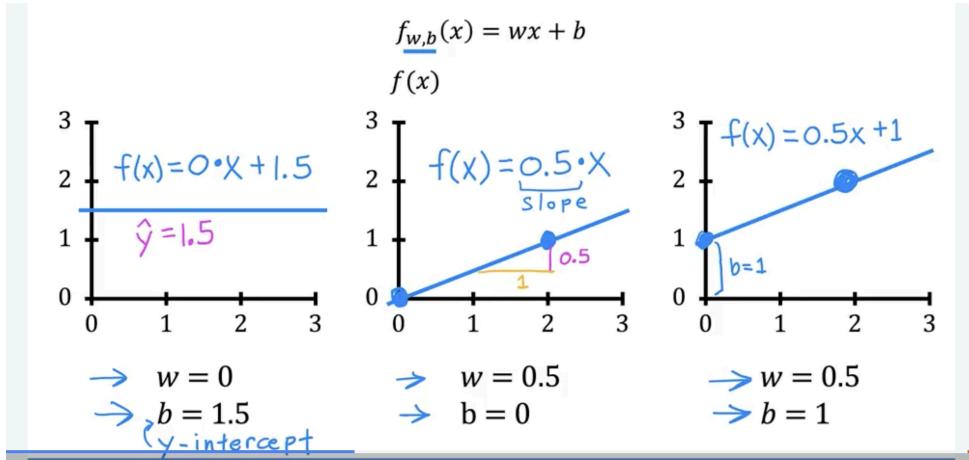
Term	Meaning
Model parameters	The variables w and b that can be adjusted during training
Weights / Coefficients	Alternate names for model parameters

◆ What w and b Do to the Line

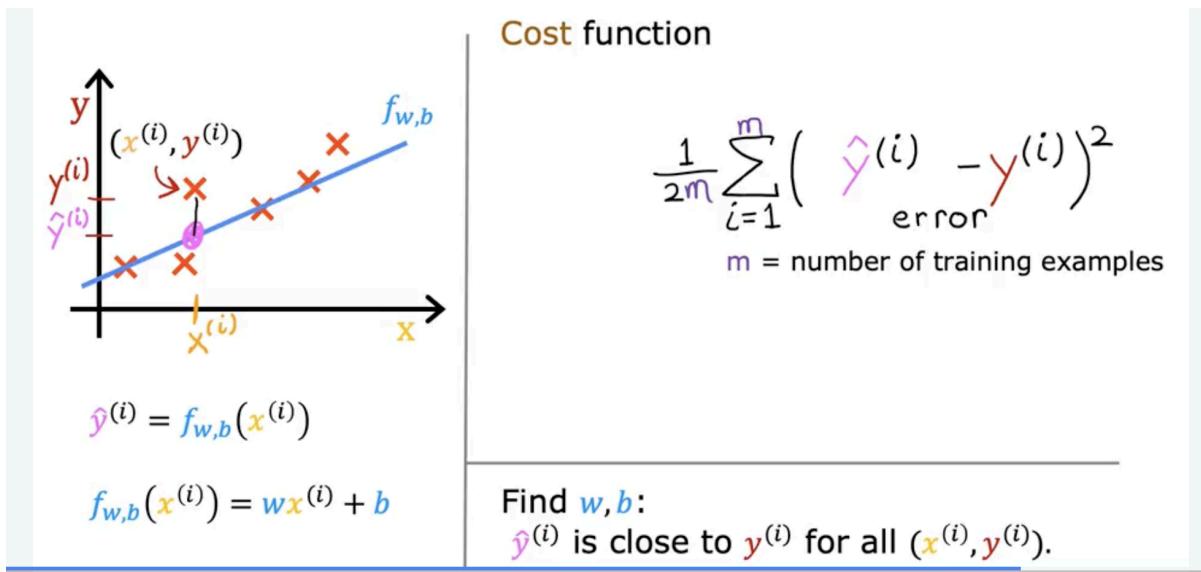
-  **Different values of w and b = different lines**
- Example 1:
 - $w = 0, b = 1.5 \rightarrow$ flat horizontal line at $y = 1.5$
 - $f(x) = 0 \cdot x + 1.5 = 1.5$
- Example 2:
 - $w = 0.5, b = 0 \rightarrow$ line through origin with slope 0.5
 - $f(x) = 0.5 \cdot x$
- Example 3:
 - $w = 0.5, b = 1 \rightarrow$ slanted line starting from $y = 1$
 - $f(x) = 0.5 \cdot x + 1$

💡 Intuition:

- w controls the **slope** of the line
- b is the **y-intercept** (where the line crosses the y-axis)



◆ Training Set Example



What's Happening?

You see:

- A blue line ($f(x) = wx + b$) – our **predicted model**
- Red Xs – the **actual training data**
- Purple points – the **predicted values \hat{y}**
- Arrows – showing **how far off the prediction is** from the real value

◆ Training Set Example

- Each training point is denoted:

$$(x^{(i)}, y^{(i)})$$

- For each input $x^{(i)}$, the model predicts:

$$\hat{y}^{(i)} = f_{w,b}(x^{(i)}) = w \cdot x^{(i)} + b$$

Goal:

Choose w and b so that predicted values $\hat{y}^{(i)}$ are as close as possible to true values $y^{(i)}$

◆ Constructing the Cost Function

❓ Why Cost Function?

To measure **how well our model fits** the training data.

◆ Step-by-Step:

◆ Step-by-Step:

1. Prediction:

$$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$$

2. Error (Difference):

$$\text{error} = \hat{y}^{(i)} - y^{(i)}$$

3. Squared Error:

$$(\hat{y}^{(i)} - y^{(i)})^2$$

4. Sum Over All Training Examples:

$$\sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

where m = number of training examples

5. Take Average to normalize:

$$\frac{1}{m} \sum_{i=1}^m \left(f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

6. Divide by 2 (convention only; helps in gradient descent math):

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m \left(f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

 This is known as the **Squared Error Cost Function**.

◆ Purpose of the Cost Function

It tells us:

"How far off are our model's predictions from the true target values?"

- If $J(w, b)$ is large → predictions are bad
- If $J(w, b)$ is small → model fits data well

🎯 Final Goal

Find the best values of w and b that minimize the cost function $J(w, b)$

V2:- Understanding the Cost Function in Linear Regression

Objective

We want to:

Find the best w and b such that the line fits the training data well.

To measure “fit,” we use the **cost function**.

Cost Function Definition

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m \left(f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

Where:

- m = number of training examples
- $x(i)$ = input
- $y(i)$ = actual output
- $f(x(i))$ = predicted output

Meaning:

The cost function $J(w, b)$ computes the average squared difference between predicted and actual values.

5. Intuition with Graphs

We plot **two graphs side-by-side**:

- **Left:** Prediction line $f_w(x)$
- **Right:** Cost function curve $J(w)$

Example Dataset:

Training set:

(1,1), (2,2), (3,3)

How Cost Changes with Different w

◆ Case 1: $w = 1$

- Prediction perfectly matches the data:

$$f(1) = 1, f(2) = 2, f(3) = 3$$

- Each squared error = 0
 - Cost = 0
- On the $J(w)$ plot, we mark point:

$$J(1) = 0$$

 Best Fit

◆ Case 2: $w = 0.5$

- Predictions: $f(1)=0.5, f(2)=1, f(3)=1.5$
- Squared errors:

$$(0.5 - 1)^2 = 0.25, (1 - 2)^2 = 1, (1.5 - 3)^2 = 2.25$$

- Sum of squared errors = 3.5
- Cost:

$$J(0.5) = \frac{3.5}{6} \approx 0.58$$

► Cost increases when the line doesn't fit well.

◆ Case 3: $w = 0$

- $f(x) = 0$ for all x
- Squared errors:

$$(1)^2 + (2)^2 + (3)^2 = 1 + 4 + 9 = 14$$

- Cost:

$$J(0) = \frac{14}{6} \approx 2.33$$

⚠ Bad fit — cost is higher.

◆ Case 4: $w = -0.5$

- Downward-sloping line
- Even worse fit → cost around **5.25**

◆ 6. Big Picture

- Each value of w defines a different prediction line.
- Each line gives a different **cost**, so it becomes a **point** on the $J(w)$ curve.
- If you **plot all values of w** , you get the **cost function curve**.

🎯 Conclusion

- The goal of training is to **find the value of w (and b) that gives the smallest cost.**

- In this example, $w = 1$ gives the lowest cost.
- Linear Regression works by **minimizing the cost function** using algorithms like gradient descent.

model:

$$f_{w,b}(x) = wx + b$$

parameters:

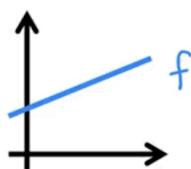
$$\underline{w, b}$$

cost function:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

goal:

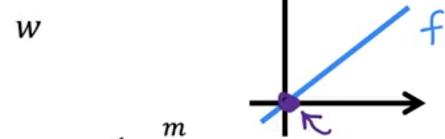
$$\underset{w, b}{\text{minimize}} J(w, b)$$



simplified

$$f_w(x) = \underline{wx}$$

$$b = \emptyset$$



$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

$$\underset{w}{\text{minimize}} J(w)$$

◆ LEFT SIDE — Full Model

♦ Model

$$f_{w,b}(x) = wx + b$$

- This is your linear regression equation
- It means: Predict output using slope w and intercept b

♦ Parameters

$$w, b$$

These are the two things you control while training:

- **w = weight (slope)**
 - **b = bias (intercept)**
-

◆ Cost Function

◆ Cost Function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m \left(f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

This measures **how bad your prediction is**:

- $f_{w,b}(x^{(i)})$ = predicted output for input $x^{(i)}$
- $y^{(i)}$ = actual output
- Square the difference, average it → you get total **cost** (or loss)

◆ Goal

$$\min_{w,b} J(w, b)$$

This means:

Find values of **w** and **b** that **minimize** the cost function.
In other words, you want your model to make **the least error possible**.

◆ RIGHT SIDE — Simplified Case

In this special case:

Assume no intercept (i.e., $b=0$)

◆ Simplified Model

$f_w(x) = wx$

- Only w matters now
 - Line always passes through origin (0, 0)
-

◆ Simplified Cost Function

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

Same as before — but now you're only adjusting w , not b .

◆ Goal

$$\min_w J(w)$$

Now you're just finding the best slope w that fits your data when the line passes through the origin.

V3:- Visualizing the Cost Function $J(w,b)$ in Linear Regression

◆ Recap: The Linear Model and Objective

- Model:

$$f_{w,b}(x) = w \cdot x + b$$

- Parameters:
 - w : slope (weight)
 - b : intercept (bias)

- Cost Function:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m \left(f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

- Goal: Find values of w and b that minimize $J(w, b)$

◆ From 1D to 2D Cost Function Visualization

Previously:

- We simplified the model by setting $b = 0$, reducing:

$$f_w(x) = w \cdot x$$

- The cost function $J(w)$ became a **U-shaped curve** (like a soup bowl) over a single variable w

Now:

- We return to the full model with both parameters:

$$f_{w,b}(x) = w \cdot x + b$$

- So, the cost function $J(w, b)$ becomes a function over two variables → 3D surface

◆ 3D Surface Visualization of $J(w,b)$

- Imagine a **soup bowl** or a **curved dinner plate** in 3D space
- Axes:
 - w : horizontal axis 1
 - b : horizontal axis 2
 - $J(w,b)$: vertical height (cost value)

⌚ Each point on this surface represents a specific combination of w and b , and its corresponding cost

Example:

- $w=-10, b=-15 \rightarrow$ height of surface = $J(-10, -15)$
-

◆ Contour Plot: Flattening the Bowl

What is a Contour Plot?

- A **2D top-down view** of the 3D cost function
- Like slicing the soup bowl **horizontally** and plotting the shape of each height
- Used in:
 - Geography (e.g. topographical maps)
 - Optimization
 - ML training visualization

Features of a Contour Plot:

Feature	Meaning
---------	---------

Concentric ovals Each oval = points with the **same** cost
 $J(w,b)J(w, b)J(w,b)$

Center oval Lowest cost (best model fit)

Further from center Higher cost (worse fits)

Analogy:

Imagine a bowl sticking **out of your screen**. If you fly directly above it and look down, you'd see **ellipses** — these are the contour lines.

◆ Visual Insights:

- Different combinations of w and b → different lines $f(x)f(x)f(x)$
 - Each line → corresponds to a **point** in the contour plot
 - Poor fit lines (bad w,bw, bw,b) = outer contour circles (higher cost)
 - Best fit line (ideal w,bw, bw,b) = center of the contours (minimum cost)
-

Using these visual tools (3D and contour plots), we understand how linear regression searches for the best parameters that give the **lowest cost** — leading to the best fit line.

V4:- Visualizing the Effect of Parameters w and b on Cost Function $J(w,b)J(w, b)J(w,b)$

◆ Connecting Line Fit to Cost $J(w,b)J(w, b)J(w,b)$

Each point on the **cost function plot** (either 3D surface or 2D contour) corresponds to:

- A specific pair of model parameters (w, b)
 - A straight line $f(x) = w \cdot x + b$
 - A value of the cost function $J(w, b)$ which tells how well that line fits the training data
-

◆ Examples of Different (w, b) Choices and Their Impact

1. Poor Fit: $w=-0.15, b=800$

Line:

$$f(x) = -0.15x + 800$$

- Effect:

- High intercept (crosses y-axis at 800)
 - Negative slope (slants downwards)
 - Far from most training points
 -  High cost $J(w, b)$
 - On the contour plot, **far from the center**
-

2. Flat Line: $w=0, b=360$

- Line:
 $f(x) = 360$
- Effect:

- Completely horizontal line
 - Ignores the value of x
 - Somewhat better than previous, but still  **not a good fit**
 - Still **far from the minimum**, but **closer** than the first
-

3. Another Poor Fit: Different w,bw, bw,b

- **Line:**
 $f(x) = wx + b$ (some random pair)
 - **Effect:**
 - Slightly better or worse than previous depending on pair
 - Still far from optimal
 - Further away from the **minimum of J**
-

4. Good Fit: Close to Minimum Cost

- **Line:**
 $f(x)$ passes close to most data points
 $f(x)$ passes close to most data points
 - **Effect:**
 - Small vertical distance (errors) between predictions and actual values
 - **Low sum of squared errors**
 -  Point is near the center of smallest contour ellipse
 - **Close to the global minimum** of cost function
-

V1:- Gradient Descent: A Beginner-Friendly Explanation

◆ Why We Need Gradient Descent

- We have a **cost function** $J(w, b)$, which tells us how good or bad our parameters w and b are.
 - Our **goal** is to find values of w and b that **minimize this cost function** — i.e., make predictions as accurate as possible.
-

◆ What is Gradient Descent?

- Gradient descent is an **optimization algorithm** used to find the **minimum** of a function.
 - It's not limited to linear regression — it's used to train deep learning models too!
-

◆ Key Idea

Think of gradient descent like:

"Walking downhill on a surface until you reach the bottom (minimum cost)."

Imagine the cost function $J(w,b)$ as a **hilly landscape**:

- You're standing at a point on this hill (initial w, b).
 - You look around in all directions.
 - You take a **tiny step downhill** in the **steepest direction** (negative gradient).
 - Repeat this process until you reach a **valley**, i.e., a **local minimum**.
-

◆ Algorithm Steps (For Linear Regression)

1. **Initialize:** Start with any value, usually
 $w=0, b=0$
 2. **Repeat** (until convergence):
 - Compute the gradient (slope) of the cost function at your current location.
 - Update the values of w and b by stepping in the opposite direction of the gradient.
 - This is called:
 $w := w - \alpha \cdot \frac{\partial J}{\partial w}$,
 $b := b - \alpha \cdot \frac{\partial J}{\partial b}$
where α is the **learning rate**.
-

◆ Local vs Global Minimum

- **Local minimum:** A valley that's minimum *in its nearby region*, but not the lowest of all.
- **Global minimum:** The **absolute lowest** point of the function.

If the cost function has multiple valleys:

- Depending on where you start (initial w,bw, bw,b), gradient descent may land in **different valleys**.
 - It doesn't jump from one valley to another — it *sticks to the one it started descending into*.
-

◆ Real-world Analogy

Think of standing on a **golf course**:

- Hills and valleys = cost surface
 - You're blindfolded and can only feel the slope around you
 - You take tiny steps downwards toward the lowest point
 - Sometimes, you might reach a small valley (local minimum), not the lowest one (global minimum)
-

◆ Important Properties

Concept	Explanation
Gradient	Direction of steepest increase
Negative Gradient	Direction of steepest decrease (used in descent)

Learning Rate (α)	Step size — too big = overshoot, too small = slow
--	--

V2:- ▼ Implementing Gradient Descent — Step-by-Step

◆ 1. Objective

We want to **minimize the cost function**

$J(w, b)$

by **iteratively updating** the parameters w and b .

◆ 2. The Update Rule

The update for each parameter looks like this:

$$w := w - \alpha \cdot \frac{d}{dw} J(w, b)$$

$$b := b - \alpha \cdot \frac{d}{db} J(w, b)$$

◆ 3. Key Terms Explained

Term	Meaning
α (Alpha)	The learning rate — controls the size of the steps. Typical values: <code>0.01</code> , <code>0.001</code> , etc.
$\frac{d}{dw} J(w, b)$	The partial derivative of the cost function w.r.t <code>w</code> — tells us the slope or direction to move.
<code>:=</code>	Assignment operator (in code: <code>=</code>). We're setting a new value to a variable.

◆ 4. Simultaneous Update (Correct Way)

This is **very important**.

✓ **Correct way (Simultaneous update):**

```
temp_w = w - alpha * dj_dw
```

```
temp_b = b - alpha * dj_db
```

```
w = temp_w
```

```
b = temp_b
```

- `dj_dw` and `dj_db` are the partial derivatives.
- We first calculate the updates **for both w and b**, then apply the changes **together**.

✗ **Incorrect way (Non-simultaneous update):**

```
temp_w = w - alpha * dj_dw
```

```
w = temp_w
```

```
temp_b = b - alpha * dj_db # This now uses the updated `w`  
(bad!)
```

```
b = temp_b
```

This version updates `w` first, and then updates `b` using the **already updated** value of `w`. That's incorrect for gradient descent.

◆ 5. Why Simultaneous Updates Matter?

Gradient descent imagines you're walking down a hill.

If you adjust one leg first and then the other — you could **fall sideways** 😞

Instead, update both **together** like taking a clean step down the slope. 🚶

◆ 6. Convergence

- We **repeat these updates** for many iterations.
 - When `w` and `b` stop changing much (i.e., **converge**) — we assume we're at a **minimum**.
-

◆ 7. Code Template (Full Loop)

Here's how you'd write this in Python:

```
# Gradient Descent Algorithm for Linear Regression
```

```
# Initialize parameters
```

```
w = 0
```

```
b = 0
```

```
alpha = 0.01 # Learning rate
```

```
# Loop for a number of iterations
```

```
for i in range(num_iterations):
```

```

# Compute derivatives

dj_dw = ... # Use actual formula

dj_db = ... # Use actual formula


# Simultaneous update

temp_w = w - alpha * dj_dw

temp_b = b - alpha * dj_db

w = temp_w

b = temp_b

```

◆ 8. Summary Table

Concept	What It Does
Cost Function $J(w, b)$	Tells us how wrong the model is
Derivative $\frac{d}{dw}, \frac{d}{db}$	Direction to update
Alpha (Learning Rate)	Step size
Simultaneous Update	Essential for correct gradient descent
Convergence	When updates become small = algorithm stops

🧠 Mental Model

Imagine you're on a hill 🏔️ trying to find the lowest point (cost).
Gradient descent tells you:

“At your current position, go slightly **left** and **a bit down**, that’s the steepest way to the bottom!”

And you **keep repeating this process** until you're no longer moving much — that's your local minimum. 

V3:- ▼ Understanding Gradient Descent Intuitively



We want to **minimize a cost function** $J(w)$ (by adjusting the parameter w) so that our model makes better predictions.

◆ Update Rule (For One Parameter)

Symbol	Meaning
w	The parameter we are trying to learn
α	Learning rate (step size)
$\frac{d}{dw} J(w)$	Derivative (slope) of the cost at current point

◆ What the Derivative Tells Us

The **derivative is the slope of the cost function** at your current point w .

- You can **visualize** it by drawing a **tangent line** at the current point on the graph of $J(w)$.
 - The **slope of that tangent** tells you which way to move to decrease the cost.
-

Case 1: Positive Derivative

Imagine:

- You start at a point where the slope is **positive**, like this: 

That means:

- $\frac{d}{dw} J(w) > 0$
- Update becomes:

$$w := w - \alpha \cdot (\text{positive number})$$

- So, w **decreases** → moves **left** on the graph.

-  This is good because the cost is **going down**.

Case 2: Negative Derivative

Now imagine:

- You start at a point where the slope is **negative**, like this: 

That means:

- $\frac{d}{dw} J(w) < 0$
- Update becomes:

$$w := w - \alpha \cdot (\text{negative number})$$

$$= w + \alpha \cdot (\text{positive number})$$

- So, w **increases** → moves **right** on the graph.

-  Again, this is good — you're getting closer to the **minimum**.

◆ Summary: What the Derivative Does

Derivative

What Happens

Positive (>0) Move **left** (decrease
 w)

Negative (<0) Move **right** (increase
 w)

This is how gradient descent **always moves toward the minimum.**

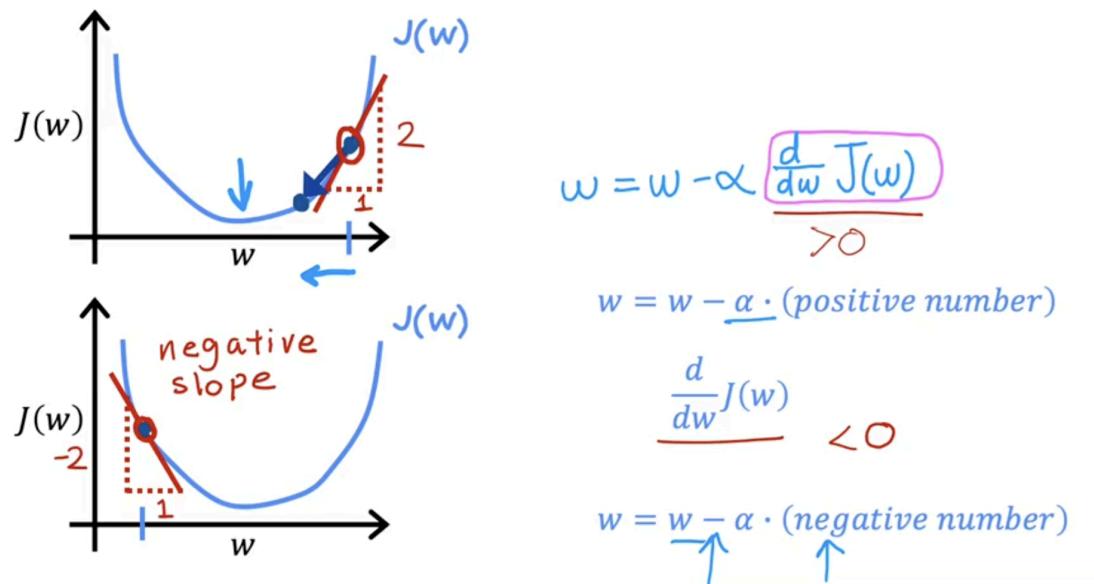
Tangent Line and Slope Intuition

When you're standing on the curve $J(w)$, the derivative is like:

"How steep is the ground under your feet?"

Gradient descent says:

"Walk in the **downhill direction**, and take a step based on how steep it is."



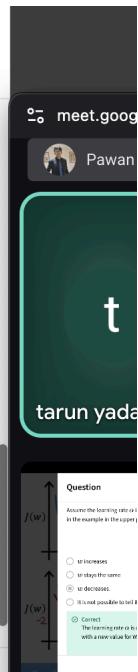
Question

Assume the learning rate α is a small positive number. When $\frac{\partial J(w,b)}{\partial w}$ is a positive number (greater than zero) -- as in the example in the upper part of the slide shown above -- what happens to w after one update step?

- w increases
- w stays the same
- w decreases.
- It is not possible to tell if w will increase or decrease.

Correct

The learning rate α is always a positive number, so if you take W minus a positive number, you end up with a new value for W that is smaller



V4:- ▶ Understanding the Learning Rate α

🧠 What is the Learning Rate?

The learning rate α controls **how big a step** gradient descent takes in the direction of the **negative gradient (downhill)**.

The update rule:

$$w := w - \alpha \cdot \nabla J(w)$$

◆ What Happens If α is Too Small?

Behavior:

- Tiny steps, very slow learning.
- The algorithm still converges, but it may take **a very long time**.
- **Inefficient** training.

Visual:

- Each update moves a **tiny bit** closer to the minimum.
- Thousands of steps may be needed.

It works but is painfully slow.

◆ What Happens If α is Too Large?

Behavior:

- Takes **huge steps**, overshoots the minimum.
- May jump back and forth across the minimum.
- Can even **diverge** (cost keeps increasing).

Visual:

- Instead of going toward the valley bottom, you leap over it repeatedly or shoot off the cliff.

It doesn't work — algorithm becomes unstable.

◆ What Happens If You're Already at the Minimum?

💡 Insight:

- At a local minimum, the **slope (derivative)** is 0.
- So, update becomes:
 $w := w - \alpha \cdot 0 \Rightarrow w := w$ —
 $w := w - \alpha \cdot 0 \Rightarrow w := w - \alpha \cdot 0 \Rightarrow w := w$
- The parameter **stops changing** — which is exactly what you want.

✓ Gradient descent naturally stops at a minimum.

◆ Why Gradient Descent Takes Smaller Steps Near the Minimum

As you move closer to a minimum:

- The **derivative** $\frac{dJ(w)}{dw}$ becomes smaller.
- So even with a **fixed α** , the **update steps get smaller** automatically.

This creates a **self-adjusting effect**:

- Large updates far from the minimum ✓
 - Smaller updates as you get closer ✓
 - Eventually, tiny updates near the bottom ✓
-

◆ Recap: Choosing α

Learning
Rate
 α

Behavior

Too small

Very slow convergence

Just right

Fast and stable convergence

Too large

Overshooting, instability, divergence

Gradient descent algorithm

Repeat until convergence

$$\left\{ \begin{array}{l} w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ b = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{array} \right.$$

Simultaneously update w and b

Assignment

$$a = c$$

$$a = a + 1$$

Code

Truth assertion

$$a = c$$

$$a = a + 1$$

Math

$$a == c$$

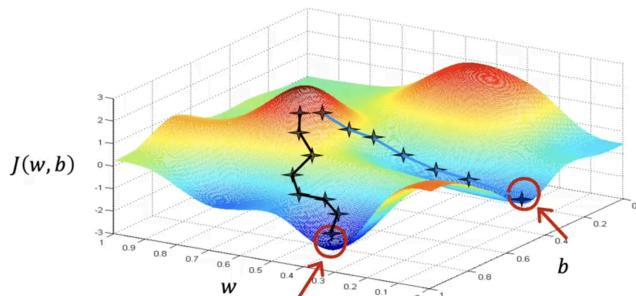
Correct: Simultaneous update

$$\left\{ \begin{array}{l} tmp_w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ tmp_b = b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w = tmp_w \\ b = tmp_b \end{array} \right.$$

Incorrect

$$\left\{ \begin{array}{l} tmp_w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ tmp_w = tmp_w \\ tmp_b = b - \alpha \frac{\partial}{\partial b} J(tmp_w, b) \\ b = tmp_b \end{array} \right.$$

More than one local minimum

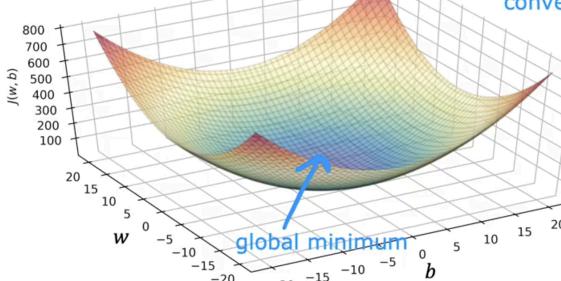


squared error cost

$J(w, b)$

bowl shape

convex function



Full Gradient Descent for Linear Regression: Step-by-Step Intuition

1. Model: Linear Regression

The model predicts output $y^{(i)}$ using:

$$\hat{y}^{(i)} = f_{w,b}(x^{(i)}) = w \cdot x^{(i)} + b$$

2. Cost Function (Mean Squared Error)

To measure how far off our predictions are:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m \left(f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

Where:

- m : number of training examples
- $f_{w,b}(x^{(i)})$: predicted value
- $y^{(i)}$: actual target value

We use $1/2m$ to simplify the derivative later.

3. Gradient Descent Algorithm

Our goal: minimize the cost $J(w,b)$
 We do this by updating w and b **iteratively**:

$$w := w - \alpha \cdot \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \cdot \frac{\partial J(w, b)}{\partial b}$$

Where:

- α : learning rate (controls step size)
-

◆ 4. Derivatives (Gradients)

Derived from calculus, but here's what you need:

$$\frac{\partial J(w, b)}{\partial w} = \frac{1}{m} \sum_{i=1}^m \left(f_{w,b}(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}$$

$$\frac{\partial J(w, b)}{\partial b} = \frac{1}{m} \sum_{i=1}^m \left(f_{w,b}(x^{(i)}) - y^{(i)} \right)$$

These tell you **how to change** w and b to **decrease the cost**.

◆ 5. Final Algorithm

Repeat until convergence (i.e. parameters don't change much):

```

# Initialize parameters

w = 0

b = 0


# For number of iterations

for i in range(num_iterations):

    # Compute predictions

    predictions = w * x + b


    # Compute errors

    errors = predictions - y


    # Compute gradients

    dj_dw = (1/m) * np.dot(errors, x)

    dj_db = (1/m) * np.sum(errors)


    # Update parameters

    w = w - alpha * dj_dw

    b = b - alpha * dj_db

```

 Make sure to do **simultaneous updates** of w and b as explained earlier.

6. Why It Works: Convex Cost Function

- For **linear regression**, the cost function $J(w,b)$ is **convex** (bowl-shaped).
- That means there is **only one global minimum** — no local traps!
- So, **no matter where you start**, gradient descent (with a good learning rate) **will find the optimal solution**.

Convex = Only one lowest point 🏆

Gradient descent will always reach it (if α isn't too large)

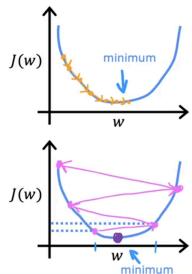
🎯 Summary: You've Just Built Your First ML Training Algorithm

Component	Meaning
Model	$y^{\wedge} = wx + b$
Cost	Measures how wrong predictions are
Gradient Descent	Algorithm to improve w , b step by step
Derivatives	Show how much and which way to change parameters
Convexity	Guarantees convergence to the global minimum

$$w = w - \alpha \frac{d}{dw} J(w)$$

If α is too small...
Gradient descent may be slow.

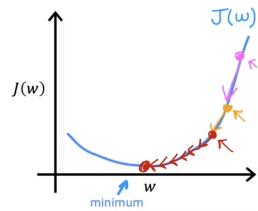
If α is too large...
Gradient descent may:
- Overshoot, never reach minimum
- Fail to converge, diverge



Can reach local minimum with fixed learning rate α

$$w = w - \alpha \frac{d}{dw} J(w)$$

smaller
not as large
large
Near a local minimum,
- Derivative becomes smaller
- Update steps become smaller
Can reach minimum without decreasing learning rate



- Gradient descent is an algorithm for finding values of parameters w and b that minimize the cost function J .

repeat until convergence {

$$\begin{aligned} w &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ b &= b - \alpha \frac{\partial}{\partial b} J(w, b) \end{aligned}$$

2. For linear regression, what is the update step for parameter b ?

When $\frac{\partial J(w, b)}{\partial w}$ is a negative number (less than zero), what happens to w after one update step?

- w increases.
- w decreases
- It is not possible to tell if w will increase or decrease.
- w stays the same

$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$

$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$

W increases